
Towards Diagram Understanding and Cognitive Reasoning in Icon Question Answering

Pan Lu¹, Liang Qiu¹, Jiaqi Chen², Tony Xia¹, Yizhou Zhao¹,
Wei Zhang³, Zhou Yu⁴, Xiaodan Liang², Song-Chun Zhu¹

¹Center for Vision, Cognition, Learning and Autonomy, UCLA

²Sun Yat-sen University, ³East China Normal University, ⁴Columbia University

Abstract

Current visual question answering (VQA) tasks mainly consider answering human-annotated questions for natural images. However, aside from natural images, abstract diagrams with semantic richness are still understudied in visual understanding and reasoning research. In this work, we introduce a new challenge of Icon Question Answering (IconQA) with the goal of answering a question in an icon image context. We release IconQA, a large-scale dataset that consists of 107,439 questions, which highlights the importance of abstract diagram understanding and comprehensive cognitive reasoning. IconQA requires not only perception skills like object recognition and text understanding, but also diverse cognitive reasoning skills, such as geometric reasoning, commonsense reasoning, and arithmetic reasoning. To facilitate potential IconQA models to learn semantic representations for icon images, we further release an icon dataset Icon645 which contains 645,687 colored icons on 377 classes. We conduct extensive user studies and blind experiments and reproduce a wide range of advanced VQA methods to benchmark the IconQA task. Also, we develop a strong IconQA baseline Patch-TRM that applies a pyramid cross-modal Transformer with input diagram embeddings pre-trained on the icon dataset. IconQA and Icon645 are available at <https://iconqa.github.io>.

1 Introduction

The long-standing goal of the VQA task is to exploit systems that can answer natural questions that correspond to visual information. Several datasets have been released to evaluate the systems' visual and textual content understanding abilities [3, 49, 11, 17, 14, 44]. One of the underlying limitations of current VQA datasets is that they are focusing on answering visual questions for natural images. However, aside from natural pictures, abstract diagrams with visual and semantic richness account for a large proportion of the visual world. Some pioneering works attempt to propose datasets that are capable of answering questions for abstract diagrams. However, these datasets either address domain-specific charts, plots, and illustrations [20, 18], or are generated from limited templates [47, 40, 17]. These limitations impede their practical applications in real-world scenarios.

To address these shortcomings, we introduce Icon Question Answering (IconQA), a new challenge for *abstract diagram* visual reasoning and question answering. The task, stemming from math word problems for children [33], exhibits a promising potential to develop education assistants. We name the proposed task as IconQA because the images depict icons, which simplify recognition and allow us to focus on reasoning skills for further research. We release IconQA, a large-scale dataset that contains 107,439 QA pairs and covers three different sub-tasks: *multiple-image-choice*, *multiple-text-choice* and *filling-in-the-blank*. A typical IconQA problem is provided with an icon image and a question, and the answer is in the form of either a short piece of text or a choice from multiple visual or textual choices. As the examples in Figure 1 show, IconQA poses new challenges

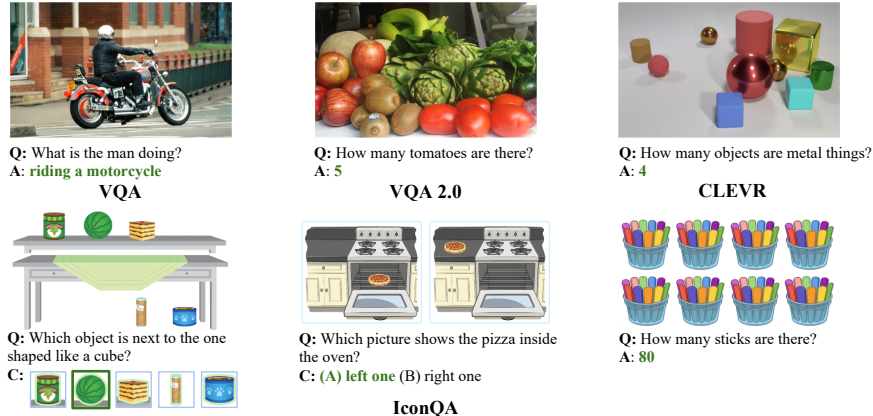


Figure 1: Examples in existing VQA datasets [3, 11, 17]) and our IconQA dataset.

for abstract diagram understanding like recognizing objects and identifying attributes. Besides, it is critical to develop diverse cognitive reasoning skills, including counting objects, comparing attributes, performing arithmetic operations, etc.

We use the IconQA dataset to benchmark various VQA approaches in the IconQA task, including four attention-based multimodal pooling methods [2, 22, 46, 9] and four Transformer-based pre-trained methods [26, 5, 45, 23]. Also, we conduct extensive user studies to evaluate the performance differences between the algorithms and human beings. Three blind studies show that the IconQA dataset is robust against biased shortcuts when answering icon questions. We further develop a strong baseline called pyramid patch cross-modal Transformer (Patch-TRM), which effectively learns implicit visual and linguistic relationships in IconQA. Along with the IconQA dataset, we collect an auxiliary icon dataset, Icon645, that features 645,687 colored icons on 377 object classes. The Icon645 dataset is used to pre-train the models to enhance abstract diagram understanding.

2 The IconQA Dataset

The IconQA dataset is collected from open-source math textbooks with rich icon images and diverse topic. We collect 107,439 IconQA data instances, where each data point contains a colored icon image, a natural language question, optional image or text choices, as well as a correct answer. The distribution of questions is visualized in Figure 2. Importantly, the diversity in the question distribution implies the requirement of high-level understanding of textual and visual contents in IconQA. Figure 3 shows the word cloud of the question text in IconQA after eliminating the stop words. The most frequent words: *shape*, *many*, and *object* indicate that answering IconQA questions requires the model to identify a variety of geometric shapes and icon objects. Inspired by this, learning informative representations for icon images plays an important role in visual reasoning for the IconQA task. See Appendix B for the details of data collection.

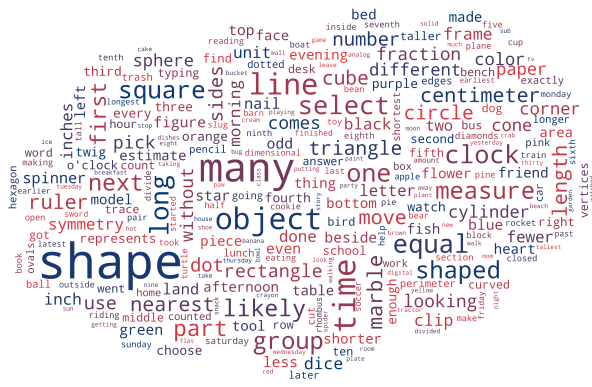
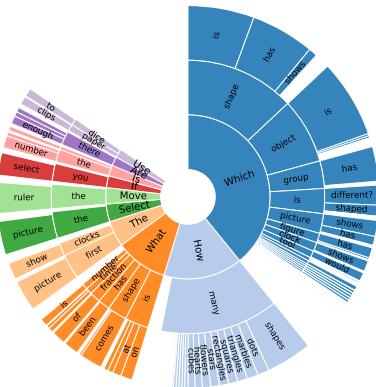


Figure 2: Question types in IconQA.

Figure 3: Word cloud of the question text in IconQA.

Skill Categories. Our IconQA dataset contains questions of multiple different cognitive reasoning and arithmetic reasoning types that can be grouped into 13 categories, shown in Table 1. We annotate

Skill types	Description
Geometry	Identify shapes, symmetry, transformations
Counting	Count objects, shapes
Comparing	Compare object attributes
Spatial	Identify spatial positions and relations
Scene	Understand abstract scenes
Pattern	Identify next and different patterns
Time	Identify time of clocks, events
Fraction	Perform fraction operations
Estimation	Estimate lengths, large numbers
Algebra	Perform algebraic operations
Measurement	Measure widths, lengths, heights
Commonsense	Apply external knowledge
Probability	Perform probability and statistics operations

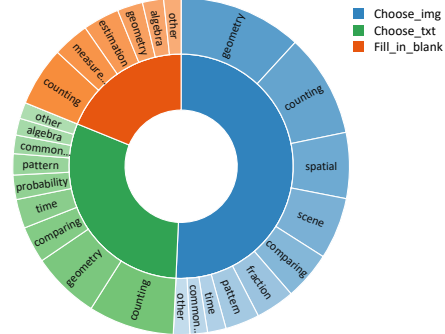


Table 1: Definition of reasoning skill types. Figure 4: Skill distribution in IconQA questions.

each question in IconQA with its corresponding skill types based on the tags provided by the original problem sources. Figure 4 shows the distributions of questions related to each skill. For instance, to answer 13.8% of the questions in IconQA, the model has to be capable of *comparing* object attributes. Additionally, each question can be related to up to three skills out of these 13 categories, and on average, a question requires 1.63 skills. The detailed statistics are demonstrated in Table 6. In general, the *filling-in-the-blank* sub-task consists of questions that require the most number of skills, averaging 1.81 skills per question. 9.25% of the *filling-in-the-blank* questions require 3 skills. Some examples that require various skills are shown in Figure 1.

Comparisons to Other Datasets. We compare our IconQA dataset with two datasets on natural images and five datasets on abstract diagrams in Table 2. To summarize, IconQA is different from these datasets in various aspects. Unlike natural images (VQA [3], CLEVR [17]) or abstract diagrams like scenes, charts, plots, and illustrations (VQA-Abstract [3], DVQA [18], NLVR [40], AI2D [20], Geometry3K [30]), IconQA features icon images and covers the largest object set of 388 classes. As questions in IconQA stem from real-world math problems and they may describe complex problem scenarios, IconQA has the longest question length among all related datasets. Furthermore, IconQA requires both commonsense and arithmetic reasoning due to its origin from real-world problems. Lastly, IconQA contains more QA task types including answering questions with image choices.

Table 2: Statistics for the IconQA dataset and comparisons with existing datasets.

	#QA	#Image	AvgQ	MaxQ	Image Type	QSource	#Object	#Task	VisualAns	CommonSen	Arithmetic
VQA [3]	614,163	204,721	6.1	23	Natural	Annotated	-	2		✓	
CLEVR [17]	999,968	100,000	18.4	43	Natural	Generated	3	1			
VQA-Abstract [3]	150,000	50,000	6.0	21	Scene	Annotated	131	2			
DVQA [18]	2,325,316	300,000	10.3	23	Bar chart	Generated	-	1			✓
NLVR [40]	92,244	92,244	11.2	25	Scatter plot	Generated	3	1			
Geometry3K [30]	3,002	2,342	10.1	46	Diagram	Real-world	4	1			✓
AI2D [20]	4,563	4,903	9.8	64	Illustration	Real-world	-	1		✓	
IconQA (Ours)	107,439	96,817	8.4	73	Icon image	Real-world	388	3	✓	✓	✓

3 Methods

Inspired by recent advances Transformer has achieved in vision-language tasks [26, 29], we develop a cross-modal Transformer model Patch-TRM for icon question answering. Taking the *multi-image choice* sub-task as an example, the overall architecture is shown in Figure 5. The diagram is first parsed into ordered patches in a hierarchical pyramid layout. These patches are then encoded by a pre-trained ResNet and passed through a vision Transformer. Question text is encoded by a language Transformer and fused with patch embeddings via the attention mechanism. The encoded image choices are concatenated with the joint diagram-question representation and then fed to a classifier for question answering. The other two sub-tasks utilize similar network architectures, except that in the *multi-text-choice* sub-task, we use an LSTM encoder [13] for choice embedding, while *filling-in-the-blank* does not need a choice encoder.

Current dominant VQA methods either rely heavily on the ResNet backbone network to extract image features or depend on the Transformer encoders to learn image embeddings. However, these networks are pre-trained on natural images and are likely to fail to extract meaningful representations or reasonable object proposals when processing the diagrams in IconQA. Instead, we pre-train the ResNet network on the icon classification task with the icon dataset we compiled (Section C). Patch-TRM hierarchically parses the diagram into patches that retain complete objects to a large

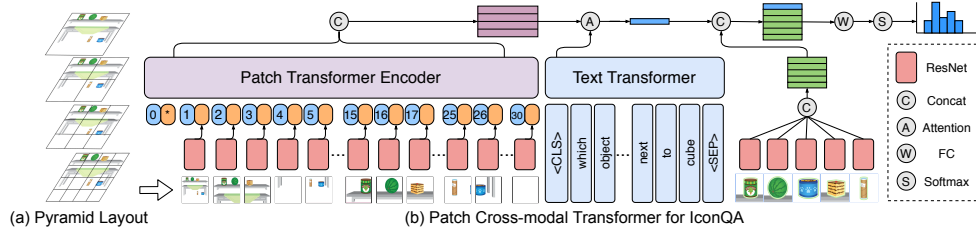


Figure 5: Our IconQA baseline Patch-TRM. Patch-TRM takes patches parsed from a hierarchical pyramid layout and embeds them through ResNet pre-trained on our Icon645 dataset. The joint diagram-question feature is learned via cross-modal Transformers followed by the attention module.

Table 3: Results on the IconQA dataset.

Method	Sub-tasks (3)			Reasoning skills (13)												
	Img.	Txt.	Blank	Geo.	Cou.	Com.	Spa.	Sc.	Pat.	Tim.	Fra.	Est.	Alg.	Mea.	Sen.	Pro.
Human	95.69	93.91	93.56	94.63	97.63	94.41	93.31	92.73	95.66	97.94	97.45	87.51	96.29	86.55	97.06	85.67
Q-Only	41.64	36.86	28.45	38.03	33.63	48.19	37.14	35.37	33.66	48.09	33.06	40.46	28.02	38.07	45.25	40.76
I-Only	41.56	36.02	46.65	38.71	37.64	45.26	37.52	35.47	36.29	47.37	32.48	62.29	31.73	64.02	45.25	37.51
Top-Down [2]	75.92	68.51	73.03	80.07	65.01	80.65	45.78	58.22	55.01	68.28	72.43	99.54	50.00	99.46	84.54	83.75
BAN [22]	76.33	70.82	75.54	79.99	67.56	82.12	53.20	66.92	55.67	66.50	73.77	97.06	47.46	96.50	82.12	82.45
ViLBERT [26]	76.66	70.47	77.08	80.05	71.05	75.60	49.46	58.52	62.78	66.72	74.09	99.22	50.62	99.07	81.78	70.94
MCAN [46]	77.36	71.25	74.52	79.86	68.94	82.73	49.70	62.49	54.79	68.00	76.20	99.08	47.32	98.99	83.25	84.87
DFAF [9]	77.72	72.17	78.28	81.80	70.68	81.69	51.42	67.01	56.60	67.72	77.60	99.02	50.27	98.83	84.11	85.70
UNITER [5]	78.71	72.39	78.53	81.31	71.01	83.67	48.34	61.25	60.81	69.77	78.37	99.41	49.18	99.38	86.10	87.84
ViT [45]	79.15	72.34	78.92	82.60	70.84	82.12	54.64	68.80	58.46	68.66	77.41	98.95	51.10	98.76	84.72	86.07
ViLT [23]	79.67	72.69	79.27	82.61	71.13	84.95	53.38	66.72	59.22	69.99	75.81	99.02	50.55	98.91	86.10	87.65
Patch-TRM (Ours)	82.66	75.19	83.62	81.87	77.81	87.00	55.62	62.39	68.75	77.98	82.13	98.24	56.73	97.98	92.49	95.73

extent, and the parsed patches are embedded by the pre-trained ResNet network before being fed into the vision Transformer. The hierarchical parsing structure, along with the ResNet pre-trained on icon data facilitate our Patch-TRM to learn informative diagram representations for the IconQA task. More details of the pre-training task are discussed in Section E.5.

4 Experiments

Table 3 demonstrates the results of the benchmark methods and our baseline on the IconQA test set. The first three columns of the results represent the three sub-tasks: *multi-image-choice*, *multi-text-choice*, and *filling-in-the-blank* respectively. The remaining 13 columns illustrate the results of these approaches over problems that require different reasoning skills, as defined in Table 1. More details of experimental settings and analysis are in Appendix E.2.

Humans outperform all benchmarks consistently over there sub-tasks and most reasoning skills. There is still a large gap to fill for future research of abstract diagram understanding and visual reasoning on the icon domain. The results achieved in blind studies of Q-only and I-only are close to random, showing that the IconQA dataset is robust and reliable in distribution. Our proposed Patch-TRM baseline outperforms current state-of-the-art VQA models in all three sub-tasks. These improvements mainly come from two insights: pre-training ResNet on icon images and taking a hierarchical approach with attention mechanism. Similarly, the Patch-TRM baseline obtains better results than the benchmarks over most reasoning skill types. Interestingly, in some skills such as *estimation*, *measurement*, and *probability*, Patch-TRM performs better than average human beings. This implies that neural networks have a promising potential to develop the basic ability of mathematical reasoning.

5 Conclusion

In this work, we introduce IconQA, an open-source dataset of icon question answering, which consists of 107,439 questions, three sub-tasks, and thirteen types of cognitive reasoning skills. We benchmark the IconQA task extensively with a user study, three blind studies, as well as multiple existing VQA approaches. We further develop a strong baseline, Patch-TRM, which parses the diagram in a pyramid layout and applies cross-modal Transformers with attention mechanism to learn the meaningful joint diagram-question feature. Additionally, we introduce Icon645, a large-scale icon dataset that is useful to pre-train the diagram encoding network used in Patch-TRM for the IconQA task.

References

- [1] Vedika Agarwal, Rakshith Shetty, and Mario Fritz. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9690–9698, 2020.
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision (CVPR)*, pages 2425–2433, 2015.
- [4] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1567–1578, 2019.
- [5] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [6] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [7] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9268–9277, 2019.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT (NAACL)*, 2018.
- [9] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6639–6648, 2019.
- [10] Peng Gao, Hongsheng Li, Shuang Li, Pan Lu, Yikang Li, Steven CH Hoi, and Xiaogang Wang. Question-guided hybrid convolution for visual question answering. In *The European Conference on Computer Vision (ECCV)*, pages 469–485, 2018.
- [11] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 770–778, 2016.
- [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [14] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019.
- [15] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. In defense of grid features for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10267–10276, 2020.
- [16] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.
- [17] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2901–2910, 2017.

- [18] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5648–5656, 2018.
- [19] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*, 2017.
- [20] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Min Joon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [21] Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4999–5007, 2017.
- [22] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1571–1581, 2018.
- [23] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 5583–5594. PMLR, 18–24 Jul 2021.
- [24] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, pages 32–73, 2017.
- [25] Li-Jia Li, Hao Su, Li Fei-Fei, and Eric P Xing. Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1378–1386, 2010.
- [26] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [27] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 2980–2988, 2017.
- [28] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2537–2546, 2019.
- [29] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13–23, 2019.
- [30] Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. In *The 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2021.
- [31] Pan Lu, Lei Ji, Wei Zhang, Nan Duan, Ming Zhou, and Jianyong Wang. R-vqa: learning visual relation facts with semantic attention for visual question answering. In *The ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pages 1880–1889, 2018.
- [32] Pan Lu, Hongsheng Li, Wei Zhang, Jianyong Wang, and Xiaogang Wang. Co-attending free-form regions and detections with multi-modal multiplicative feature embedding for visual question answering. In *The AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [33] Maria Martiniello. Language and the performance of english-language learners in math word problems. *Harvard Educational Review*, 78(2):333–368, 2008.
- [34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems (NeurIPS)*, pages 91–99, 2015.

- [35] Mrinmaya Sachan, Kumar Dubey, and Eric Xing. From textbooks to knowledge: A case study in harvesting axiomatic knowledge from textbooks to solve geometry problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 773–784, 2017.
- [36] Mrinmaya Sachan, Kumar Avinava Dubey, Tom M Mitchell, Dan Roth, and Eric P Xing. Learning pipelines with limited data and domain knowledge: A study in parsing physics problems. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 140–151, 2018.
- [37] Mike Schuster and Kaisuke Nakajima. Japanese and korean voice search. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5149–5152. IEEE, 2012.
- [38] Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. Solving geometry problems: Combining text and diagram interpretation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1466–1476, 2015.
- [39] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8317–8326, 2019.
- [40] Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 217–223, 2017.
- [41] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*, 2019.
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems (NeurIPS)*, pages 5998–6008, 2017.
- [43] Shuo Wang, Yizhou Wang, and Song-Chun Zhu. Learning hierarchical space tiling for scene modeling, parsing and attribute tagging. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(12):2478–2491, 2015.
- [44] Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [45] Ildoo Kim Wonjae Kim, Bokyung Son. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. under review.
- [46] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6281–6290, 2019.
- [47] Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and Yang: Balancing and answering binary visual questions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [48] Jun Zhu, Tianfu Wu, Song-Chun Zhu, Xiaokang Yang, and Wenjun Zhang. A reconfigurable tangram model for scene representation and categorization. *IEEE Transactions on Image Processing*, 25(1):150–166, 2015.
- [49] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

Appendix

A Related Works

VQA Datasets. There have been efforts to develop datasets for the visual question answering (VQA) task since the first large-scale benchmark was introduced in [3]. Early released datasets [11, 24, 39, 44] contain natural images and related questions, where understanding the visual and textual contents is essential for question answering. Some recent datasets introduce questions that involve more diverse visual scenes or require external knowledge to answer, which leads to more complex visual and semantic reasoning for question answering. For example, CLEVR [17] is a synthetic dataset that serves as a diagnostic test for a range of visual reasoning abilities over combinations of three object shapes. However, these datasets are limited to the natural image domain and pay little attention to abstract diagrams, which also have informative semantics and wide applications.

Diagram QA Datasets. To address the need for vision-and-language reasoning for diagrams, several abstract diagram QA datasets have been developed. For example, abstract VQA [3, 47] considers the task of answering questions on abstract scenes. Similarly, NLVR [40], FigureQA [19], and DVQA [18] feature diagrams of scientific plots that are generated with several figure types and question templates. However, questions and diagrams in these datasets are generated from limited templates, leading to the existence of unintended linguistic shortcuts for question answering. Some more works have proposed datasets of middle school math or science problems in more practical and complex scenarios [38, 21, 35, 36, 30]. A central limitation of the subject QA datasets is that they require complex domain-specific knowledge, which makes disentangling visual reasoning and domain knowledge difficult. Herein, we address these limitations by introducing the IconQA dataset, where only elementary commonsense is required. Through IconQA, we aim to provide a new benchmark for abstract scene understanding and learning different visual reasoning skills in *real-world* scenarios.

VQA Methods. Early approaches on visual question answering usually combine multi-modal inputs by applying attention mechanisms over image regions or question words [22, 32, 31, 10, 46, 9]. Inspired by the semantic nature of VQA images, a line of approaches adopt object proposals from pre-trained object detectors and learn their semantic relationships [22, 46, 9]. As Transformers achieve excellent performance on vision tasks, pioneering works have attempted to use pre-trained models to learn visual representations for natural images in the VQA task [29, 26, 5, 23] and achieve significant improvements. However, current VQA models are not capable of extracting meaningful visual representations from abstract diagrams, as they require image embeddings or object proposals learned from natural images. Instead, we develop a strong baseline that feeds spatial patch sequences into a Transformer encoder that is powered by the embedding module pre-trained on our Icon645 dataset.

B The IconQA Dataset

The IconQA dataset provides diverse questions that require abstract diagram recognition, comprehensive visual reasoning skills, and basic commonsense knowledge. IconQA consists of 107,439 questions split across three different sub-tasks. To the best of our knowledge, IconQA is the largest VQA dataset that focuses on real-world problems with icon images while involving multiple human intelligence reasoning abilities.

B.1 Data Collection

We aim to collect icon-based question answering pairs that involve multiple reasoning skills, such as visual reasoning and commonsense reasoning. To construct the IconQA dataset, which stems from real-world math word problems, we search for open-source math textbooks with rich icon images and diverse topics. Of those, we choose *IXL Math Learning* which compiles popular textbooks aligned to California Common Core Content Standards¹. We ask well-trained crowd workers to collect problems that cover content from pre-K to third grade, as these problems usually contain abstract images and involve little to none complex domain knowledge. With the driven interest of visual reasoning over abstract images, we filter out the questions that do not accompany icon images or only have images

¹<https://www.ixl.com/standards/california/math>

in black and white. Redundant or repetitive data instances are also removed. Question choices are randomly shuffled to ensure a balanced answer distribution.

The IconQA dataset consists of 107,439 questions and is divided into train, validation, and test splits with a ratio of 6:2:2, as shown in Table 5. The dataset consists of three sub-tasks: *multi-image-choice*, *multi-text-choice*, and *filling-in-the-blank*. The *multi-image-choice* sub-task is defined as choosing the correct image from a list of image candidates based on a given diagram and its corresponding question. Similarly, the *multi-text-choice* sub-task is defined as a multiple choice question with 2-5 text choices and an abstract diagram. The *filling-in-the-blank* sub-task is similar to the common VQA task, requiring a brief text answer for each question, except in IconQA, the images are icon images instead of natural images.

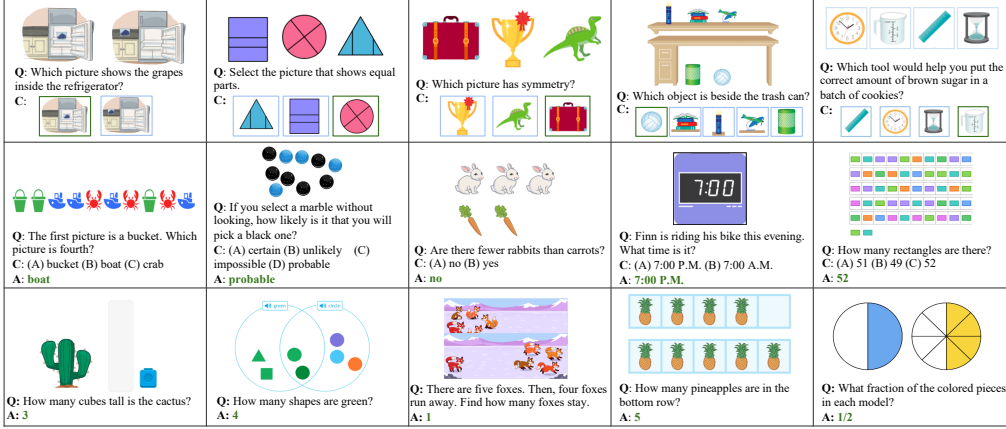


Figure 6: More examples in the IconQA dataset. **Top:** The *multi-image-choice* sub-task. **Middle:** The *multi-text-choice* sub-task. **Bottom:** The *filling-in-the-blank* sub-task.

Question Skill Categories. The questions we collected contain meta-information including question topics, chapter names, image names, etc. After extensive data exploration by well-informed individuals, we designed a set of rules that map each question to 1-3 of the 13 categories based on trigger words in metadata. The rules for trigger words are list in Table 4.

Table 4: Trigger words in metadata for skill categories.

Skill types	Trigger words in metadata
Geometry	name the shape, shapes of, classify shapes, solid, corners, faces, edges, vertices, sides, dimensional, rectangle, circle, triangle, square, rhombus, sphere, cylinder, cone, cubes, hexagon, perimeter, area, curved, open and close, flip turn, symmetry
Counting	count, tally, a group, ordinal number, area, even or odd, place value, represent numbers, comparing review, equal sides, square corners, one more, one less, fewer, enough, more.
Comparing	compare, comparing, more, less, fewer, enough, wide and narrow, light and heavy, long and short, tall and short, match analog and digital
Spatial	top, above, below, beside, next to, inside and outside, left
Scene	problems with pictures, beside, above, inside and outside, wide and narrow, objects
Pattern	the next, comes next, ordinal number, different
Time	clock, am or pm, elapsed time, times
Fraction	equal parts, halves, thirds, fourths, fraction
Estimation	estimate, measure
Algebra	count to fill, skip count, tally, even or odd, tens and ones, thousands, of ten, elapsed time, perimeter, area, divide
Measurement	measure
Commonsense	light and heavy, compare size, holds more or less, am or pm, times of, tool
Probability	likely

B.2 Dataset Analysis

Figure 7 (a) illustrates the distribution of question lengths of each sub-task in the IconQA dataset. For simplicity, all questions longer than 35 words are counted as having 35 words. Questions in

the *multi-text-choice* sub-task distribute more evenly, while for *multi-img-choice*, there is a long-tail distribution due to the complexity of textual scenarios. We find that some icon objects are frequently mentioned in the questions. In Figure 7 (b), the frequencies of the 40 most frequently mentioned icons are shown. These icon entities cover different daily-life objects such as animals, plants, shapes, food, etc. We cluster question sentences into different types based on frequent trigram prefixes starting the sentences.

Table 5: Statistics for the IconQA dataset.

Tasks	All	Train	Val	Test
<i>Multi-image-choice</i>	57,672	34,603	11,535	11,535
<i>Multi-text-choice</i>	31,578	18,946	6,316	6,316
<i>Filling-in-the-blank</i>	18,189	10,913	3,638	3,638
All	107,439	64,462	21,489	21,489

Table 6: Skill numbers for questions in IconQA.

Task	Avg.	1 skill	2 skills	3 skills
<i>Multi-image-choice</i>	1.51	55.78%	37.44%	6.77%
<i>Multi-text-choice</i>	1.73	33.21%	60.14%	6.65%
<i>Filling-in-the-blank</i>	1.81	28.30%	62.43%	9.25%
All	1.63	44.50%	48.34%	7.16%

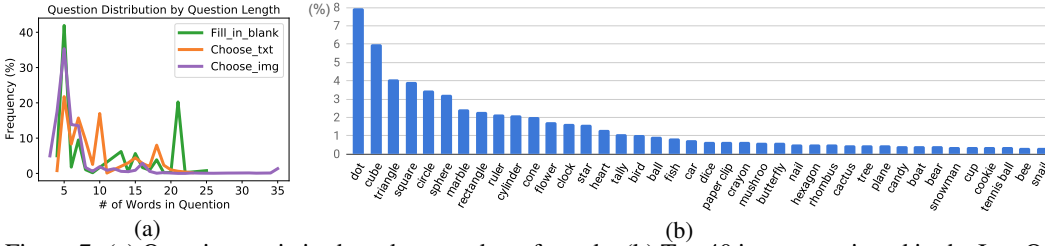


Figure 7: (a) Question statistics based on number of words. (b) Top 40 icons mentioned in the IconQA question texts and their appearance percentage. These icons cover various types of real-world objects.

IconQA presents new challenges in icon understanding and cognitive reasoning to many existing visual reasoning methods. 1) Icons feature intrinsic natures of abstract symbolism, varied styles, and ambiguous semantics, which differs from natural images significantly. 2) Since there is a lack of high-quality annotation data for icon diagrams, it restricts current mainstream data-driven visual methods to transfer smoothly to the icon domain. 3) As 107,439 questions in IconQA stem from real-world math word problems, it has made 13 different cognitive reasoning skills essential, including spatial reasoning, commonsense reasoning, estimation, and arithmetic calculation.

B.3 User Study

Using Amazon Mechanical Turk (AMT), we ask people to provide answers to the questions in the test set along with their age group. We also strongly encourage parents who have young children to let their children complete the questionnaires, as their answers give us insights to how the designed audience of these questions perform. The test set is split into batches of 20 questions, which we call a task, with each task assigned to 3 crowd workers on AMT. This amounts to a total of 64,467 effective test set answers.

To ensure the truthfulness of the age information, we ask the participants to select their age at both the beginning and the end of the questionnaire, with the age choices appearing in 2 different orders. To ensure the quality of the answers, we include 4 attention check questions: 3 of which are about the instructions, making sure that the participants read the instructions carefully. We also add an extra fake question in the middle for each *choosing an image choice* and *choosing a text choice* task, instructing them to choose the fourth choice despite what the choices are. Figure 8 shows the instructions and the first three attention check questions. Figure 9 shows the fake question along with the age confirmation. Figure 10, 11, and 12 are example questions for three sub-tasks respectively. We also make sure that the workers answering our tasks have a history HIT approval rate of at least 95% and a previous approval count of 1,000.

In summary, for each Human Intelligence Task (HIT) on AMT, we have 2 age questions, 4 attention check questions, and 20 real questions from the IconQA test set. Among the 64,467 test answers, we filter out 1) the questionnaires that do not pass the 4 attention check questions, 2) the questionnaires that do not answer consistently for the two age-related questions, 3) the questionnaires that are finished unreasonably slowly/quickly. After filtering, we have 54,896 effective question answers, which we believe is a decently large sample for the human performance study.

Overview

Thank you for helping us with our research!

- You will be answering **21 multiple choice questions** in the following task within **20 minutes**.
- For each question, there will be **1 image** providing some context information, and there will be **2-6 image choices** to select from.
- Please refer to the image and try your best to pick the **one best answer** with the information from the image.
- If a particular question seems ambiguous (no correct answer/more than one correct answer/etc.), please choose the answer that makes the most sense to you.
- We will be collecting your age group information purely for research purposes. Be assured that the information will be stored anonymously and will not be tied to you.
- If you **have young children**, we encourage you to let them try and **answer all the questions individually** in the HIT. It would help us greatly.
- Please select the last buttons in the following three questions to proceed.

I have read the Overview carefully and will answer to my best capability.

- ☐ No
☐ Yes

How many questions will be in this questionnaire?

- ☐ 6
☐ 11
☐ 16
☐ 21

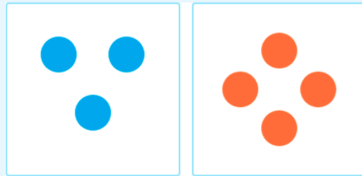
If you have a child, we strongly encourage you to

- ☐ Work together with your child.
☐ Let your child finish the tasks individually.

Submit

Figure 8: AMT instructions for the user study.

12. Please select the fourth choice in the following question.



Choices and your answer:



☐ choice 1



☐ choice 2



☐ choice 3



☐ choice 4



☐ choice 5

Just to confirm, what's your (child's) age? If you are letting your child answer the questions, please specify the child's age.

- ☐ 9 - 18
☐ 3 - 8
☐ 19+

Figure 9: AMT attention check questions.

Instruction: Given an image, select the image choice that best answers the question.

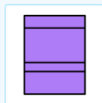
4. Select the picture that shows equal parts.



Choices and your answer:



☐ choice 1



☐ choice 2

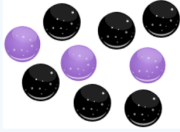


☐ choice 3

Figure 10: An AMT question example for the *multi-image-choice* sub-task.

Instruction: Given an image, select the text choice that best answers the question.

20. If you select a marble without looking, how likely is it that you will pick a black one?



Choices and your answer:

☐ probable

☐ certain

☐ impossible

☐ unlikely

Figure 11: An AMT question example for the *multi-text-choice* sub-task.

Instruction: Given an image, give your concise answer to the question.

15. How many triangles are there?



Your Answer











Figure 12: An AMT question example for the *filling-in-the-blank* sub-task.

C The Icon645 Dataset

As discussed in Section B.2, IconQA questions are accompanied by abstract diagrams that cover a wide range of icon objects. Using existing backbone networks to extract image representations for these icon images is inadequate, as most of these networks are pre-trained on natural images. To overcome the limitation, we develop a new large-scale icon dataset for pre-training existing vision backbone networks. We use the collected icon data to pre-train the current backbone networks, which can be applied to extract diagram representations in IconQA.

We retrieve the 388 icon classes mentioned in the question texts from FlatIcon², the largest database of free vector icons. After removing 11 classes that can't be retrieved, we construct an icon dataset containing 377 classes, called Icon645. The Icon645 dataset includes 645,687 colored icons with a minimum size of 64 by 64 and a maximum size of 256 by 256. Examples in Table 7 show that our collected icons include a wide variety of colors, formats and styles. On top of pre-training encoders, the large-scale icon data could also contribute to future research on abstract aesthetics and symbolic visual understanding. In this work, we use the icon data to pre-train backbone networks on the icon classification task in order to extract semantic representations from abstract diagrams in IconQA.

Table 7: Collected icon examples in the Icon645 dataset.

Icons	Examples	Icons	Examples
Bed		Bucket	
Cake		Car	
Castle		Dog	
Giraffe		Kite	
Soda		Tree	

D Details of Patch-TRM

We develop a patch cross-modal Transformer model (Patch-TRM) as a strong baseline for the IconQA task as illustrated in Figure 5. We will introduce the details of Patch-TRM as follows.

²FlatIcon: <https://www.flaticon.com/>

D.1 Diagram Encoder

Similar to natural images in most VQA datasets, abstract diagrams also have rich visual and semantic information that is critical to answering questions. Current dominant VQA methods [3, 2, 22, 9, 46, 15, 1] either extract high-level visual representations from a pre-trained ResNet backbone network [12] in a top-down fashion, or apply a bottom-up mechanism to extract semantic representations via an object detector, such as a model based on Faster R-CNN [34]. However, these methods depend heavily on the backbone network, which is pre-trained on natural images. When processing diagrams in IconQA, they are likely to fail to extract meaningful representations or reasonable object proposals. Inspired by the early progress in using hierarchical scene layout to parse images [25, 48, 43] and the recent advances in Transformer-based image encoding [29, 26, 45], we develop a method that splits diagrams into hierarchical patch sequences from a pyramid structure and learns their visual representations using a visual Transformer.

As diagrams in IconQA have more varied aspect ratios than natural images, we add blank paddings at the bottom or on the right side of the images to ensure that they are square-shaped. Each padded diagram is then cropped into a set of patch sequences with different scales. The padding operation and the hierarchical scene layout can facilitate extracting complete objects that retain specific semantics. Let $p = [p_1, p_2, \dots, p_n]$ denote the patch sequence in the splitting order from the original diagram. From each patch sequence, we extract the visual features using a ResNet model and represent the features as $f_p = [f_{p_1}, f_{p_2}, \dots, f_{p_n}]$. The representation for each patch, f_{p_i} , is then summed up with its positional embedding with respect to its sequential index i . Finally, the updated visual patch embeddings pass through a standard multi-layer Transformer [42] to learn high-level visual representations $h_p = [h_{[CLS]}, h_{p_1}, h_{p_2}, \dots, h_{p_n}]$. Here, the trainable token [CLS], which is added to the Transformer inputs, learns the global meaning of these sequences. As mentioned before, it is not feasible to use existing pre-trained ResNet to process abstract diagrams due to a lack of similar resources for pre-training. So we pre-train the ResNet on icon classification with the icon dataset we compiled (Section C). More details of the pre-training task are discussed in Section E.5.

D.2 Language Encoder

Questions in IconQA have a wide distribution of question lengths, so we follow the recent approaches [42, 16, 41, 26, 29] that apply the BERT model [8] to embed question texts, rather than using traditional LSTM [13] or GRU [6] for long sequence encoding. Given the question w_0, w_1, \dots, w_t , the input is formatted as $[[CLS], w_0, w_1, \dots, w_t]$. We use the WordPiece [37] subword tokenizer and the resulting sequence is padded to the maximum length. Similar to other methods that use BERT as sentence encoders, we consider the output corresponding to the first token [CLS] as the embedding of the entire question, noted as h_q .

D.3 Answer Reasoning

Given the image patch representation $h_p \in \mathcal{R}^{n \times k}$, and question embedding $h_q \in \mathcal{R}^k$, where n denotes the number of diagram patches and k denotes the learned embedding size of the patches, we apply a cross-modal attention to learn their joint representation:

$$a = \text{softmax}(W_p h_p \circ W_q h_q), \quad (1)$$

$$h_v = \sum_i^n a(i) \times h_{p_i}, \quad (2)$$

where W_p and W_q are learnable mapping parameters, and \circ is the element-wise product operator. The joint representation h_v is calculated as the weighted sum over all diagram patches.

Before predicating the answer, multiple choice candidates need to be encoded. Taking the *multi-image-choice* task as an example, each image choice is encoded as the output of the last pooling layer of the pre-trained ResNet. The encoded image choice is denoted as $h_c \in \mathcal{R}^{m \times k}$, where m is the number of the candidates. The choice embeddings are concatenated with the diagram-question representation, and then the resulted embeddings are fed to a classifier over the candidates:

$$p_{ans} = \text{softmax}(W_a([h_v; h_c]) + b_a), \quad (3)$$

where W_a and b_a are classifier parameters, and p_{ans} is the probability of the predicated answer choice.

Similarly, in the *multi-text-choice* sub-task, the answer is predicated over text choices, except that each text choice is embedded with LSTM layers first. We formulate the *filling-in-the-blank* sub-task as a multi-class classification problem from all possible answers in the training data, as most VQA works do. After generating the joint encoding for the input diagram and question, a linear classifier is trained to predict the final answer.

E Experiments

E.1 Benchmark Methods

In this section, we first develop a patch cross-modal Transformer model (Patch-TRM) as a strong baseline for the IconQA task. To benchmark the IconQA dataset, we consider multi-modal pooling methods with attention mechanisms [2, 22, 9, 46], Transformer-based VQA approaches [29, 5, 45, 23], and three blind study methods as benchmark models, as summarized in Figure 13. Additionally, a user study is conducted to explore the performances of human beings in different age groups. In the sections below, we discuss the main principles of the core networks in the benchmarks we performed.

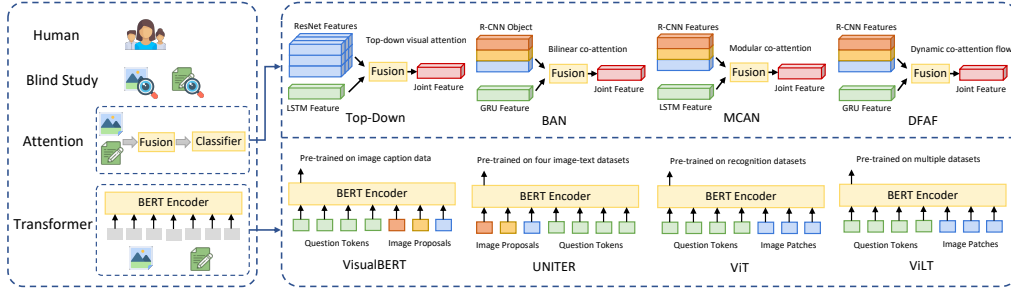


Figure 13: An overview of benchmark baselines on the IconQA task.

Attention models. We construct four attention models for benchmarking. The first model implements Top-Down attention [2] for VQA, which is a strong attention method that applies free-form based attention on image representations from a pre-trained ResNet-101 network. The remaining three models utilize the bottom-up attention mechanism with the help of object detection proposals from Faster-RCNN [34]. Specifically, BAN [22] proposes a method that utilizes bilinear attention distributions to learn joint vision-language information. DFAF [9] is an advanced model that applies self-attention and cross-modal attention and introduces the information flow to help the model focus on target question words and image regions. The last approach, MCAN [46], learns the self-attention on the questions and images and the question-guided-attention of images jointly.

Transformer models. Four Transformer-based models are also implemented as benchmarks. ViLBERT [29] and UNITER [5] are two Transformer-based approaches that take image object proposals from Faster-RCNN [34] and question tokens as inputs. Specifically, ViLBERT learns the joint representation of the image content and the natural language content from image proposal regions and question tokens, while UNITER processes multimodal inputs simultaneously for joint visual and textual understanding. The last two benchmarks ViL [45] and ViLT [23] are more recently proposed Transformer models that take image patch tokens instead of object proposals as inputs when representing the image.

Blind study models. We develop three models to check for possible data biases in the IconQA dataset. A random baseline picks up one from the given choice candidates for the *multiple-choice* sub-tasks while predicts the answer by randomly selecting one from all possible answers in the train data for the *filling-in-the-blank* sub-task. Q-Only is set up similar to the Top-Down [2] model, but it only considers textual inputs. This baseline learns the question bias in the training set. I-Only also has a Top-Down architecture, but it only takes abstract diagrams as inputs, and tests the distribution biases in the images and answers in IconQA.

E.2 Experimental Details

Following prior work [3], all the baselines are trained on the IconQA training set, tuned on the validation set, and finally evaluated on the test set. Similar to [3], we choose accuracy as the evaluation metric. For the two *multi-choice* sub-tasks, the answer is regarded as correct only if it matches the ground truth. On the other hand, as the collected answers for *filling-in-blank* are short numbers, correct answers are expanded to include both the digital number and its corresponding words. More details of the benchmark setups and implementations can be found in Appendix E.2.

Our benchmarks and baselines are implemented using PyTorch. All experiments are run on one Nvidia RTX 3090 GPU. We use the Adamax optimizer with optimal learning rates of 7×10^{-4} , 8×10^{-4} , and 2×10^{-3} on the three sub-tasks respectively. We apply a binary cross-entropy loss to train the multi-class classifier with a batch size of 64 and a maximum epoch of 50. The early stopping strategy is used when the validation accuracy stops improving for five consecutive epochs. It takes about 50, 30, and 10 minutes to train our baseline Patch-TRM on three sub-tasks respectively. We use the same learning parameters set in Top-Down [2] when evaluating the eight benchmarks and our developed baseline Patch-TRM. Some crucial parameters used in our model are clarified below.

Our Baseline Model. For our baseline Patch-TRM, each diagram is split four times by varied scales, resulting in 79 (1+4+9+16+49) patches totally. After resizing them to 224×224 , patch visual features are extracted from the last pooling layer, resulting in a 2048-d feature vector. The ResNet network used to embed the patches is pre-trained on the icon classification task as discussed in Section E.5. The patch Transformer has one layer of Transformer block with four attention heads and outputs embeddings with a hidden state size of 768. A small pre-trained BERT model [41] is used to encode the question text in the language encoder.

Attention models. For Top-Down, the attention-based baselines use $7 \times 7 \times 2048$ -d features from the last convolution layer. For BAN [22], DFAF [9], and MCAN [46], image features of dimension 2,048 are extracted from Faster R-CNN [34]. Question words in these attention models are encoded into features of dimension 1,024 by GRU [6]. And the visual and textual features are then embedded into 1,024 dimensions with the corresponding attention mechanisms and fusion methods reported in original works.

Transformer models. For ViLBERT [29] and UNITER [5], we use Faster R-CNN [34] to extract 36 proposal regions as the visual inputs. Both ViL [45] and ViLT [23] use ViT-B/32 pre-trained on ImageNet to encode the image embeddings. The hidden size is set as 768, the layer depth is 32, and the input image is sliced into patches with a size of 32. For ViL, we use two dependent Transformers to embed the question and image respectively.

E.3 Experiment Results

Human Performance. Out of the 54,896 collected answers, 9,620 are made by young children from age 3 to 8, 19,040 are made by adolescents from age 9 to 18, and 26,236 are made by adults. The human performance over the three sub-tasks and thirteen skills is illustrated in Figure 14. As expected, young children do not answer the questions as well as adolescents or adults, suggesting that most participants answer their ages correctly. Moreover, the result shows that humans perform more consistently on all sub-tasks compared to machine algorithms. Interestingly, humans are outperformed by models quite significantly in questions that require numerical reasoning skills like *probability*, *measurement*, and *estimation*.

E.4 Ablation Study

To study the functions of individual components in our model, we conduct an ablation analysis. Table 8 presents the results of different simplifications of our full model, where each implementation is trained on the IconQA train set and tested on the validation set. Instead of ResNet101 pre-trained on the icon classification task, *Patch-TRM w/o pre* utilizes ResNet101 pre-trained on natural image data for patch feature extraction. The decreasing performance of 0.95-2.49% indicates that pre-training backbones on tasks within similar domains is critical to downstream tasks. The attention mechanism helps to combine the image and question representations and improves the model performance by up to 7% compared to using simple concatenation (denoted as *Patch-TRM w/o att*). Positional embeddings of the ordered diagram patches benefit the vision Transformer by enabling it to learn

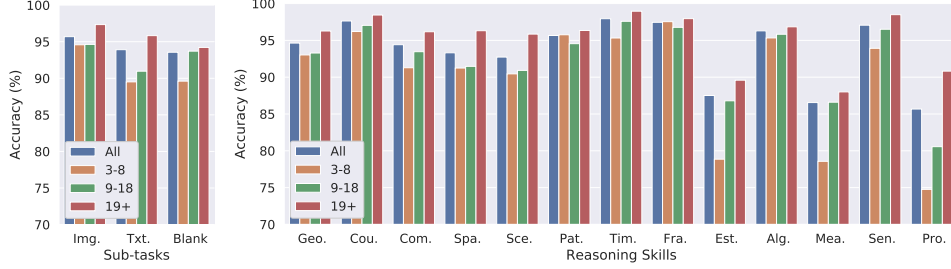


Figure 14: Performance of humans in different age groups for the IconQA task. **Left:** Accuracy over three sub-tasks; **Right:** Accuracy over thirteen reasoning skills.

spatial relationships among the patches, compared to the baseline without position embeddings (*Patch-TRM w/o pos*). *Patch-TRM V-CLS* uses the output embedding of [CLS] token as the diagram feature instead, which leads to a drastic performance decline. We have also experimented with coarse-grained patch cropping (e.g., *Pyramid 1+4+9+16* denotes 30 patches, *Pyramid 1+4+9* denotes 14 patches), which results in a performance degradation of 0.51% to 7.79%.

Q: Which object is next to the one shaped like a cube?

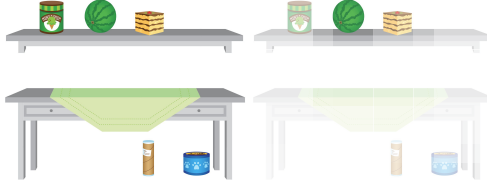


Figure 15: Text-to-image attention visualization.

Method	Img.	Txt.	Blank
Patch-TRM w/o pre	82.01	72.72	81.67
Patch-TRM w/o att	80.63	68.00	80.29
Patch-TRM w/o pos	81.27	64.98	80.68
Patch-TRM V-CLS	80.15	63.90	70.27
Pyramid 1+4+9+16	82.45	68.76	82.19
Pyramid 1+4+9	80.61	67.42	81.36
Full model	82.96	75.21	83.10

Table 8: Ablation study in IconQA.

E.5 Icon Classification for Pre-training

The Icon645 dataset is collected to pre-train the backbone network for patch feature extraction. The dataset has a long-tailed distribution, and thus we address the class-imbalanced issue following previous studies on specific loss functions such as CB loss [7], Focal loss [27], and LDAM loss [4]. The metric of Top-5 accuracy is used to evaluate different model setups and the evaluation results are summarized in Table 9. Following [28], to demonstrate performances on different data parts, we divide the dataset into three balanced clusters: Head, Medium, and Tail, corresponding to 132, 122, and 123 classes respectively. All classes in Head have at least 1,000 instances, all classes in Medium have 300 - 999 instances, and all classes in Tail have fewer than 300 instances. As the results show, the backbone network ResNet101 with a re-balanced LDAM loss function achieves the best result for icon classification on Icon645. Consequently, we adopt this pre-trained ResNet101 network to extract patch features in our baseline Patch-TRM for IconQA.

Table 9: Results for icon classification.

Method	Total	Head	Medium	Tail
ResNet32 [12] + CB [7]	27.91	19.66	36.51	33.53
ResNet32 [12] + Focal Loss [27]	32.80	51.59	36.51	8.94
ResNet32 [12] + LDAM [4]	42.65	55.68	46.42	24.94
ResNet101 [12] + LDAM [4]	62.93	70.29	70.50	47.51

E.6 Quantitative Analysis

We visualize one example with the cross-modal attention map generated by our baseline Patch-TRM in Figure 15. The visualized attention shows that our baseline is capable of attending to the corresponding patch regions with higher weights given the input question.

Figure 16 presents five examples from the IconQA test set predicted by our Patch-TRM baseline for each sub-task. Although Patch-TRM achieves promising results for most problems in IconQA, it still fails to address some complicated cases. For example, it encounters difficulties in identifying dense objects and making multi-hop reasoning.





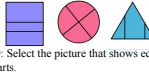











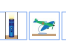
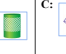
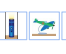
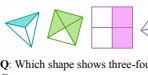









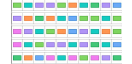
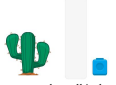


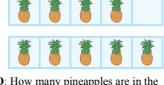

 Q: Which picture shows the grapes inside the refrigerator? C:   Ours: 	 Q: Select the picture that shows equal parts. C:    Ours: 	 Q: Which picture has symmetry? C:    Ours: 	 Q: Which object is beside the trash can? C:    Ours: 	 Q: Which shape shows three-fourths? C:     Ours: 
 Q: The first picture is a bucket. Which picture is fourth? C: (A) bucket (B) boat (C) crab Ours: <i>boat</i>	 Q: If you select a marble without looking, how likely is it that you will pick a black one? C: (A) certain (B) unlikely (C) impossible (D) probable Ours: <i>probable</i>	 Q: Are there fewer rabbits than carrots? C: (A) no (B) yes Ours: <i>no</i>	 Q: Finn is riding his bike this evening. What time is it? C: (A) 7:00 P.M. (B) 7:00 A.M. Ours: <i>7:00 P.M.</i>	 Q: How many rectangles are there? C: (A) 51 (B) 49 (C) 52 Ours: <i>51</i>
 Q: How many cubes tall is the cactus? Ours: <i>3</i>	 Q: How many shapes are green? Ours: <i>4</i>	 Q: How many faces does this shape have? Ours: <i>6</i>	 Q: How many pineapples are in the bottom row? Ours: <i>5</i>	 Q: How many blocks are there? Ours: <i>16</i>

Figure 16: Result examples predicted by our Patch-TRM model in the IconQA test set. **Top:** *Multi-image-choice* sub-task. **Middle:** *Multi-text-choice* sub-task. **Bottom:** *Filling-in-the-blank* sub-task. Correctly predicted answers are highlighted by green, while wrong ones are highlighted by red.